

Configuration Manual

MSc Research Project
Cloud Computing

Mohammed Zubair Shaik
X23228946

School of Computing
National College of Ireland

**National College of Ireland
Project Submission Sheet
School of Computing**



Student Name	Mohammed Zubair Shaik
Student ID	23228946
Programme	MSc in Cloud Computing
Year	2024-2025
Module	MSc in Research Computing
Supervisor	Yasantha Samarawickrama
Submission Due Date	11/08/2025
Project Title	Cloud-Based Intrusion Detection using Super Learner Ensemble in ICS/SCADA
Word Count	712
Page Count	16

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Mohammed Zuabir
Date:	11/08/2025

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Cloud-Based Intrusion Detection using Super Learner

Ensemble in ICS/SCADA

Mohammed Zubair Shaik

Student ID: x23228946

1. Introduction

This project develops a machine learning-based Intrusion Detection System (IDS) for ICS/SCADA networks using the ToN-IoT dataset, which includes telemetry, audit logs, and network traffic data under normal and attack scenarios. The pipeline includes data preprocessing, visualization, model training (Random Forest, Gradient Boosting, Naïve Bayes, Logistic Regression, Super Learner), evaluation, and optional cloud deployment using AWS services.

2. Dataset Information

Dataset Link: <https://research.unsw.edu.au/projects/toniot-datasets>

he datasets can be used for validating and testing various Cybersecurity applications-based AI such as intrusion detection systems, threat intelligence, malware detection, fraud detection, privacy-preservation, digital forensics, adversarial machine learning, and threat hunting.

The datasets have been called 'ToN_IoT' as they include heterogeneous data sources collected from Telemetry datasets of IoT and IIoT sensors, Operating systems datasets of Windows 7 and 10 as well as Ubuntu 14 and 18 TLS and Network traffic datasets. The datasets were collected from a realistic and large-scale network designed at the Cyber Range and IoT Labs, the School of Engineering and Information technology (SEIT), UNSW Canberra @ the Australian Defence Force Academy (ADFA). A new testbed network was created for the industry 4.0 network that includes IoT and IIoT networks. The testbed was deployed using multiple virtual machines and hosts of windows, Linux and Kali operating systems to manage the interconnection between the three layers of IoT, Cloud and Edge/Fog systems. Various attacking techniques, such as DoS, DDoS and ransomware, against web applications, IoT gateways and computer systems across the IoT/IIoT network. The datasets were gathered in a parallel processing to collect several normal and cyber-attack events from network traffic, Windows audit traces, Linux audit traces, and telemetry data of IoT services.

3. System Configuration

3.1 Software Specification

a. Python==3.10.12:

Python is an interpreted , high-level programming language . Its high-level built in datastructures , combined with dynamic typing and dynamic binding , make it very attractive for Rapid Application Development , as well as for use as a scripting or glue language to connect existing components together. Its language concept and object-oriented approach aim to help programmers write clear, logical code for small and large-scale projects

b) Visual Studio Code:-

Visual Studio Code (famously known as VS Code) is a free open source text editor by Microsoft. VS Code is available for Windows, Linux, and macOS.(Visual Studio Code -Code Editing. Redefined)

c) Google Colab:-

Collaboratory, or “Colab” for short, is a product from Google Research. Colab allows anybody to write and execute arbitrary python code through the browser, and is especially well suited to machine learning, data analysis and education. More technically, Colab is a hosted Jupyter notebook service that requires no setup to use, while providing access free of charge to computing resources including GPUs.(colab.google,)

3.2 Hardware Specification

1. Operating system: Windows 11
2. Processor: > 11th Gen Intel(R) Core(TM) i3-1115G4 @ 3.00GHz 3.00 GHz
3. System Compatibility: 64-bit
4. Hard Disk: 512GB
5. RAM: 16 GB

4. Libraries Configuration

1. Pandas
2. Numpy
3. Matplotlib
4. Seaborn
5. Plotly
6. Scikit-learn
7. Tensorflow
9. Keras

5. Project Implementation

5.1 Cloud Environment Setup

- **Cloud9:** Configured as the main development IDE.
- **EC2:** Deployed for running and testing the implementation.
- **S3 Bucket:** Created and configured for storing datasets and models.

5.2 Data Collection

- 150 CSV files were collected and stored in the S3 bucket.
- Loaded into the Cloud9 environment using boto3 and pandas.

5.3 Data Preparation

- Cleaned data by handling missing values.
- Removed irrelevant columns.
- Ensured data consistency.

5.4 Data Preprocessing

- Handled missing values in the dataset.
- Encoded categorical variables for model compatibility.
- Normalized numerical features for consistent scaling.
- Removed irrelevant or redundant columns.

5.5 Model Training

Implemented and trained multiple models:

- Random Forest
- Gradient Boosting
- Naïve Bayes
- Logistic Regression
- Super Learner ensemble

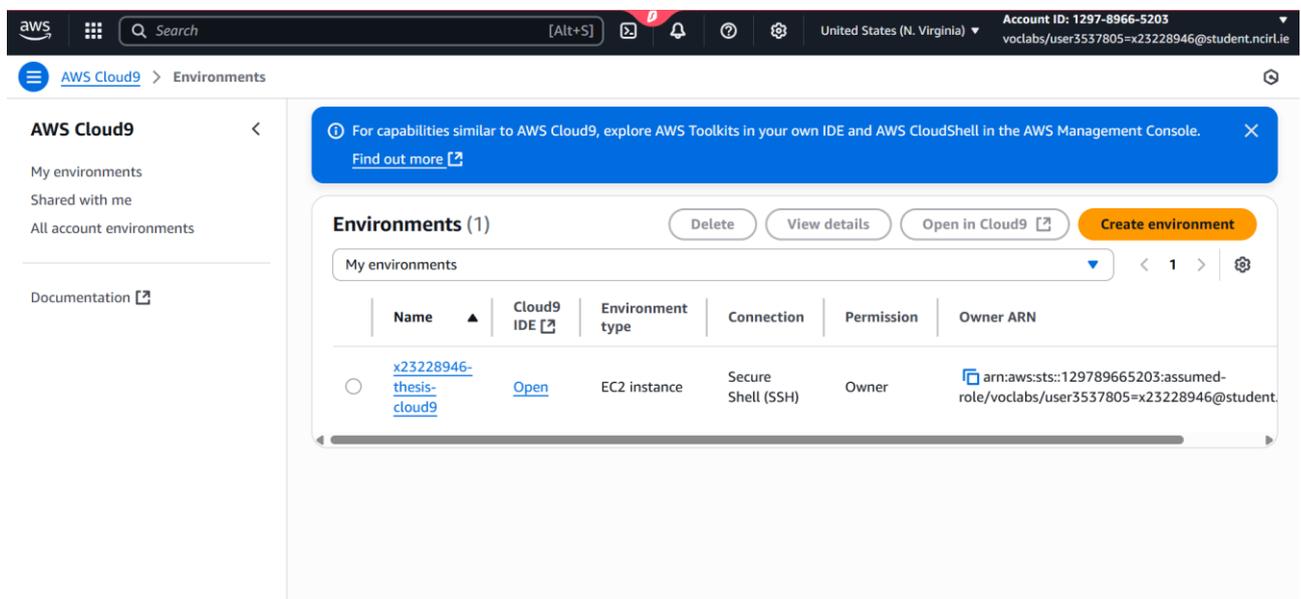
Split dataset into 80% training and 20% testing sets.

5.6 Performance Evaluation

- Metrics Used: Accuracy, Precision, Recall, F1-score, ROC-AUC.
- Generated confusion matrices and classification reports for each model.

6. Visual References

Figure 1: Cloud9 Setup



Cloud9 Enviromen

```

1  {
2  "cells": [
3  {
4  "cell_type": "code",
5  "execution_count": 1,
6  "metadata": {},
7  "outputs": [
8  {
9  "name": "stdout",
10 "output_type": "stream",
11 "text": [
12 "Defaulting to user installation because normal site-packages is not writeable\n",
13 "Requirement already satisfied: boto3 in /home/ec2-user/.local/lib/python3.9/site-packages (1.39.3)\n",
14 "Requirement already satisfied: pandas in /home/ec2-user/.local/lib/python3.9/site-packages (2.3.0)\n",
15 "Requirement already satisfied: seaborn in /home/ec2-user/.local/lib/python3.9/site-packages (0.13.2)\n",
16 "Requirement already satisfied: plotly in /home/ec2-user/.local/lib/python3.9/site-packages (6.2.0)\n",
17 "Requirement already satisfied: imblearn in /home/ec2-user/.local/lib/python3.9/site-packages (0.8)\n",
18 "Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in /usr/lib/python3.9/site-packages (from boto3) (0.10.0)\n",
19 ]
20 }
21 ]
22 }
23 ]
24 }
25 ]
26 }
27 ]
28 }
29 ]
30 }
31 ]
32 }
33 ]
34 }
35 ]
36 }
37 ]
38 }
39 ]
40 }
41 ]
42 }
43 ]
44 }
45 ]
46 }
47 ]
48 }
49 ]
50 }
51 ]
52 }
53 ]
54 }
55 ]
56 }
57 ]
58 }
59 ]
60 }
61 ]
62 }
63 ]
64 }
65 ]
66 }
67 ]
68 }
69 ]
70 }
71 ]
72 }
73 ]
74 }
75 ]
76 }
77 ]
78 }
79 ]
80 }
81 ]
82 }
83 ]
84 }
85 ]
86 }
87 ]
88 }
89 ]
90 }
91 ]
92 }
93 ]
94 }
95 ]
96 }
97 ]
98 }
99 ]
100 }
101 ]
102 }
103 ]
104 }
105 ]
106 }
107 ]
108 }
109 ]
110 }
111 ]
112 }
113 ]
114 }
115 ]
116 }
117 ]
118 }
119 ]
120 }
121 ]
122 }
123 ]
124 }
125 ]
126 }
127 ]
128 }
129 ]
130 }
131 ]
132 }
133 ]
134 }
135 ]
136 }
137 ]
138 }
139 ]
140 }
141 ]
142 }
143 ]
144 }
145 ]
146 }
147 ]
148 }
149 ]
150 }
151 ]
152 }
153 ]
154 }
155 ]
156 }
157 ]
158 }
159 ]
160 }
161 ]
162 }
163 ]
164 }
165 ]
166 }
167 ]
168 }
169 ]
170 }
171 ]
172 }
173 ]
174 }
175 ]
176 }
177 ]
178 }
179 ]
180 }
181 ]
182 }
183 ]
184 }
185 ]
186 }
187 ]
188 }
189 ]
190 }
191 ]
192 }
193 ]
194 }
195 ]
196 }
197 ]
198 }
199 ]
200 }
201 ]
202 }
203 ]
204 }
205 ]
206 }
207 ]
208 }
209 ]
210 }
211 ]
212 }
213 ]
214 }
215 ]
216 }
217 ]
218 }
219 ]
220 }
221 ]
222 }
223 ]
224 }
225 ]
226 }
227 ]
228 }
229 ]
230 }
231 ]
232 }
233 ]
234 }
235 ]
236 }
237 ]
238 }
239 ]
240 }
241 ]
242 }
243 ]
244 }
245 ]
246 }
247 ]
248 }
249 ]
250 }
251 ]
252 }
253 ]
254 }
255 ]
256 }
257 ]
258 }
259 ]
260 }
261 ]
262 }
263 ]
264 }
265 ]
266 }
267 ]
268 }
269 ]
270 }
271 ]
272 }
273 ]
274 }
275 ]
276 }
277 ]
278 }
279 ]
280 }
281 ]
282 }
283 ]
284 }
285 ]
286 }
287 ]
288 }
289 ]
290 }
291 ]
292 }
293 ]
294 }
295 ]
296 }
297 ]
298 }
299 ]
300 }
301 ]
302 }
303 ]
304 }
305 ]
306 }
307 ]
308 }
309 ]
310 }
311 ]
312 }
313 ]
314 }
315 ]
316 }
317 ]
318 }
319 ]
320 }
321 ]
322 }
323 ]
324 }
325 ]
326 }
327 ]
328 }
329 ]
330 }
331 ]
332 }
333 ]
334 }
335 ]
336 }
337 ]
338 }
339 ]
340 }
341 ]
342 }
343 ]
344 }
345 ]
346 }
347 ]
348 }
349 ]
350 }
351 ]
352 }
353 ]
354 }
355 ]
356 }
357 ]
358 }
359 ]
360 }
361 ]
362 }
363 ]
364 }
365 ]
366 }
367 ]
368 }
369 ]
370 }
371 ]
372 }
373 ]
374 }
375 ]
376 }
377 ]
378 }
379 ]
380 }
381 ]
382 }
383 ]
384 }
385 ]
386 }
387 ]
388 }
389 ]
390 }
391 ]
392 }
393 ]
394 }
395 ]
396 }
397 ]
398 }
399 ]
400 }
401 ]
402 }
403 ]
404 }
405 ]
406 }
407 ]
408 }
409 ]
410 }
411 ]
412 }
413 ]
414 }
415 ]
416 }
417 ]
418 }
419 ]
420 }
421 ]
422 }
423 ]
424 }
425 ]
426 }
427 ]
428 }
429 ]
430 }
431 ]
432 }
433 ]
434 }
435 ]
436 }
437 ]
438 }
439 ]
440 }
441 ]
442 }
443 ]
444 }
445 ]
446 }
447 ]
448 }
449 ]
450 }
451 ]
452 }
453 ]
454 }
455 ]
456 }
457 ]
458 }
459 ]
460 }
461 ]
462 }
463 ]
464 }
465 ]
466 }
467 ]
468 }
469 ]
470 }
471 ]
472 }
473 ]
474 }
475 ]
476 }
477 ]
478 }
479 ]
480 }
481 ]
482 }
483 ]
484 }
485 ]
486 }
487 ]
488 }
489 ]
490 }
491 ]
492 }
493 ]
494 }
495 ]
496 }
497 ]
498 }
499 ]
500 }
501 ]
502 }
503 ]
504 }
505 ]
506 }
507 ]
508 }
509 ]
510 }
511 ]
512 }
513 ]
514 }
515 ]
516 }
517 ]
518 }
519 ]
520 }
521 ]
522 }
523 ]
524 }
525 ]
526 }
527 ]
528 }
529 ]
530 }
531 ]
532 }
533 ]
534 }
535 ]
536 }
537 ]
538 }
539 ]
540 }
541 ]
542 }
543 ]
544 }
545 ]
546 }
547 ]
548 }
549 ]
550 }
551 ]
552 }
553 ]
554 }
555 ]
556 }
557 ]
558 }
559 ]
560 }
561 ]
562 }
563 ]
564 }
565 ]
566 }
567 ]
568 }
569 ]
570 }
571 ]
572 }
573 ]
574 }
575 ]
576 }
577 ]
578 }
579 ]
580 }
581 ]
582 }
583 ]
584 }
585 ]
586 }
587 ]
588 }
589 ]
590 }
591 ]
592 }
593 ]
594 }
595 ]
596 }
597 ]
598 }
599 ]
600 }
601 ]
602 }
603 ]
604 }
605 ]
606 }
607 ]
608 }
609 ]
610 }
611 ]
612 }
613 ]
614 }
615 ]
616 }
617 ]
618 }
619 ]
620 }
621 ]
622 }
623 ]
624 }
625 ]
626 }
627 ]
628 }
629 ]
630 }
631 ]
632 }
633 ]
634 }
635 ]
636 }
637 ]
638 }
639 ]
640 }
641 ]
642 }
643 ]
644 }
645 ]
646 }
647 ]
648 }
649 ]
650 }
651 ]
652 }
653 ]
654 }
655 ]
656 }
657 ]
658 }
659 ]
660 }
661 ]
662 }
663 ]
664 }
665 ]
666 }
667 ]
668 }
669 ]
670 }
671 ]
672 }
673 ]
674 }
675 ]
676 }
677 ]
678 }
679 ]
680 }
681 ]
682 }
683 ]
684 }
685 ]
686 }
687 ]
688 }
689 ]
690 }
691 ]
692 }
693 ]
694 }
695 ]
696 }
697 ]
698 }
699 ]
700 }
701 ]
702 }
703 ]
704 }
705 ]
706 }
707 ]
708 }
709 ]
710 }
711 ]
712 }
713 ]
714 }
715 ]
716 }
717 ]
718 }
719 ]
720 }
721 ]
722 }
723 ]
724 }
725 ]
726 }
727 ]
728 }
729 ]
730 }
731 ]
732 }
733 ]
734 }
735 ]
736 }
737 ]
738 }
739 ]
740 }
741 ]
742 }
743 ]
744 }
745 ]
746 }
747 ]
748 }
749 ]
750 }
751 ]
752 }
753 ]
754 }
755 ]
756 }
757 ]
758 }
759 ]
760 }
761 ]
762 }
763 ]
764 }
765 ]
766 }
767 ]
768 }
769 ]
770 }
771 ]
772 }
773 ]
774 }
775 ]
776 }
777 ]
778 }
779 ]
780 }
781 ]
782 }
783 ]
784 }
785 ]
786 }
787 ]
788 }
789 ]
790 }
791 ]
792 }
793 ]
794 }
795 ]
796 }
797 ]
798 }
799 ]
800 }
801 ]
802 }
803 ]
804 }
805 ]
806 }
807 ]
808 }
809 ]
810 }
811 ]
812 }
813 ]
814 }
815 ]
816 }
817 ]
818 }
819 ]
820 }
821 ]
822 }
823 ]
824 }
825 ]
826 }
827 ]
828 }
829 ]
830 }
831 ]
832 }
833 ]
834 }
835 ]
836 }
837 ]
838 }
839 ]
840 }
841 ]
842 }
843 ]
844 }
845 ]
846 }
847 ]
848 }
849 ]
850 }
851 ]
852 }
853 ]
854 }
855 ]
856 }
857 ]
858 }
859 ]
860 }
861 ]
862 }
863 ]
864 }
865 ]
866 }
867 ]
868 }
869 ]
870 }
871 ]
872 }
873 ]
874 }
875 ]
876 }
877 ]
878 }
879 ]
880 }
881 ]
882 }
883 ]
884 }
885 ]
886 }
887 ]
888 }
889 ]
890 }
891 ]
892 }
893 ]
894 }
895 ]
896 }
897 ]
898 }
899 ]
900 }
901 ]
902 }
903 ]
904 }
905 ]
906 }
907 ]
908 }
909 ]
910 }
911 ]
912 }
913 ]
914 }
915 ]
916 }
917 ]
918 }
919 ]
920 }
921 ]
922 }
923 ]
924 }
925 ]
926 }
927 ]
928 }
929 ]
930 }
931 ]
932 }
933 ]
934 }
935 ]
936 }
937 ]
938 }
939 ]
940 }
941 ]
942 }
943 ]
944 }
945 ]
946 }
947 ]
948 }
949 ]
950 }
951 ]
952 }
953 ]
954 }
955 ]
956 }
957 ]
958 }
959 ]
960 }
961 ]
962 }
963 ]
964 }
965 ]
966 }
967 ]
968 }
969 ]
970 }
971 ]
972 }
973 ]
974 }
975 ]
976 }
977 ]
978 }
979 ]
980 }
981 ]
982 }
983 ]
984 }
985 ]
986 }
987 ]
988 }
989 ]
990 }
991 ]
992 }
993 ]
994 }
995 ]
996 }
997 ]
998 }
999 ]
1000 }

```

Now we need to use URL for accessing the jupyter Notebook

Instead of 127.0.0.1 we need to use EC2 Public IPv4 address i.e., 3.238.106.93

```

python3 - "ip-172-31-64-71" x bash - "ip-172-31-64-78" ec x

To access the server, open this file in a browser:
file:///home/ec2-user/.local/share/jupyter/runtime/jpserver-3028-open.html
Or copy and paste one of these URLs:
http://ip-172-31-64-78.ec2.internal:8080/tree?token=21e310a2e3d0f1bcbcf40e7378bc0adbfc3f6d56c11607
http://127.0.0.1:8080/tree?token=21e310a2e3d0f1bcbcf40e7378bc0adbfc3f6d56c11607
[I 2025-08-10 21:25:36.165 ServerApp] Skipped non-installed server(s): bash-language-server, dockerfile-language-server-nodejs, javascript-typescript-languageserver, jedi-language-server, julia-language-server, pyright, python-language-server, python-lsp-server, r-languageserver, sql-language-server, texlab, typescript-language-server, unified-language-server, vscode-css-languageserver-bin, vscode-html-languageserver-bin, vscode-json-languageserver-bin, yamll-language-server
[I 2025-08-10 21:25:43.801 ServerApp] 302 GET / (@5.181.2.18) 0.77ms
[I 2025-08-10 21:25:43.928 JupyterNotebookApp] 302 GET /tree? (@5.181.2.18) 0.94ms

```

Figure 2: S3 Bucket Configuration

us-east-1.console.aws.amazon.com/s3/buckets/x23228946-thesis-bucket?region=us-east-1&tab=objects&bucketType=general

Account ID: 1297-8966-5203
voclabs/user3537805=x23228946@student.ncirl.ie

Amazon S3

- General purpose buckets
- Directory buckets
- Table buckets
- Vector buckets [Preview](#)
- Access Grants
- Access Points (General Purpose Buckets, FSx file systems)
- Access Points (Directory Buckets)
- Object Lambda Access Points
- Multi-Region Access Points
- Batch Operations
- IAM Access Analyzer for S3

Block Public Access settings for this account

Storage Lens

CloudShell Feedback

© 2025, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

General purpose buckets (1) [Info](#)

Buckets are containers for data stored in S3.

Find buckets by name

Copy ARN Empty Delete **Create bucket**

Name	AWS Region	Creation date
x23228946-thesis-bucket	US East (N. Virginia) us-east-1	July 5, 2025, 12:39:53 (UTC+01:00)

Account snapshot [Info](#)

Updated daily

Storage Lens provides visibility into storage usage and activity trends.

[View dashboard](#)

External access summary - new [Info](#)

Updated daily

External access findings help you identify bucket permissions that allow public access or access from other AWS accounts.

us-east-1.console.aws.amazon.com/s3/buckets/x23228946-thesis-bucket?region=us-east-1&tab=objects&bucketType=general

Account ID: 1297-8966-5203
voclabs/user3537805=x23228946@student.ncirl.ie

Amazon S3 > Buckets > x23228946-thesis-bucket

x23228946-thesis-bucket [Info](#)

General purpose buckets

Directory buckets

Table buckets

Vector buckets [Preview](#)

Access Grants

Access Points (General Purpose Buckets, FSx file systems)

Access Points (Directory Buckets)

Object Lambda Access Points

Multi-Region Access Points

Batch Operations

IAM Access Analyzer for S3

Block Public Access settings for this account

Storage Lens

Objects (2)

Copy S3 URI Copy URL Download Open Delete **Actions**

Create folder Upload

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

Name	Type	Last modified	Size	Storage class
graphs/	Folder	-	-	-
Train_Test_Network.csv	csv	July 5, 2025, 12:40:30 (UTC+01:00)	66.6 MB	Standard

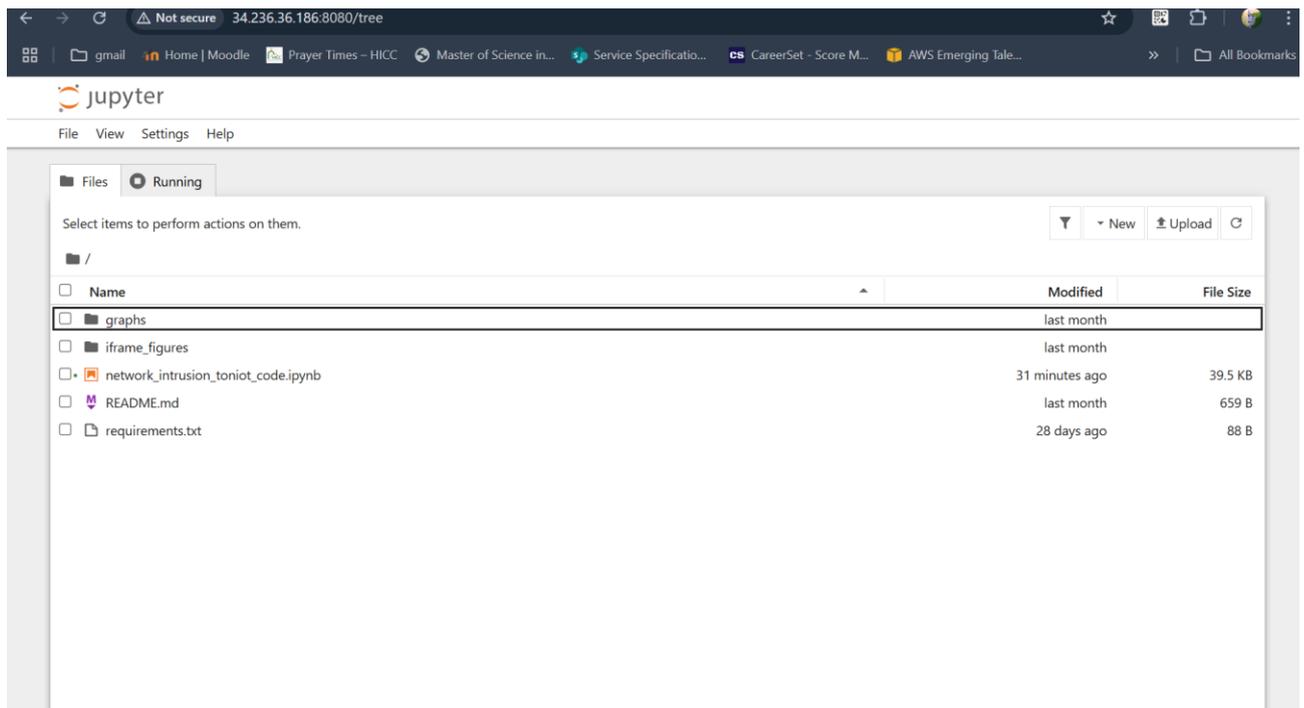
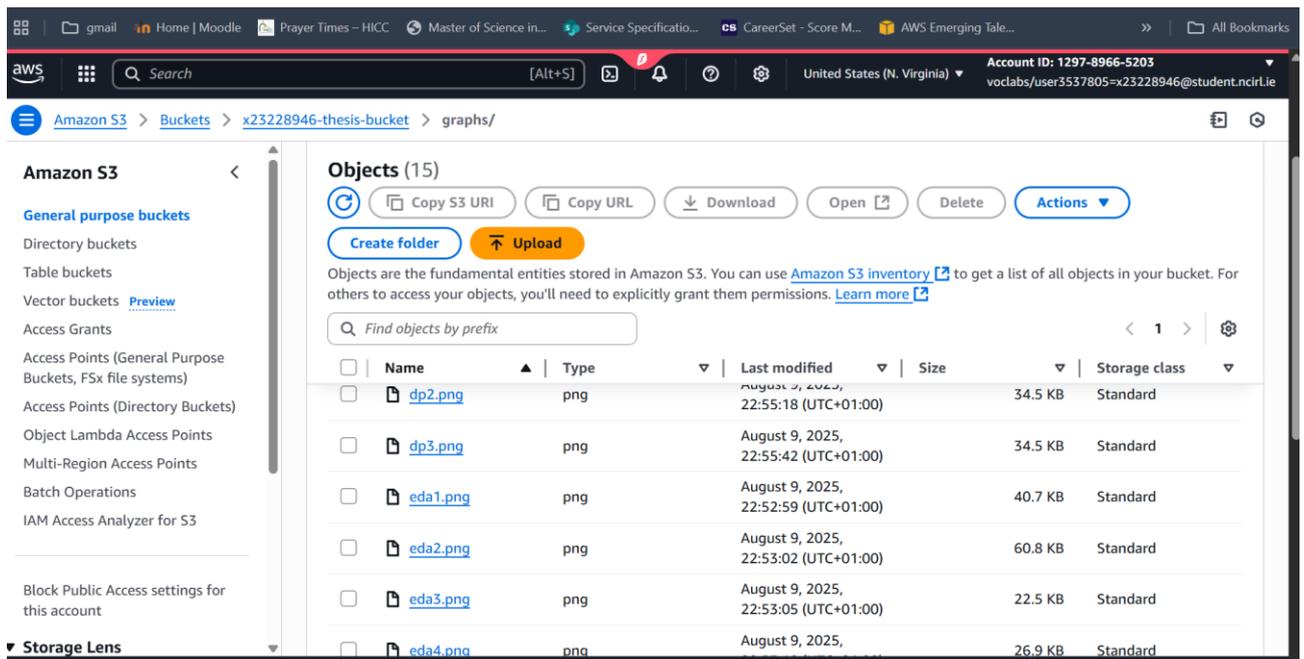


Figure 3: Ec2 Configuration

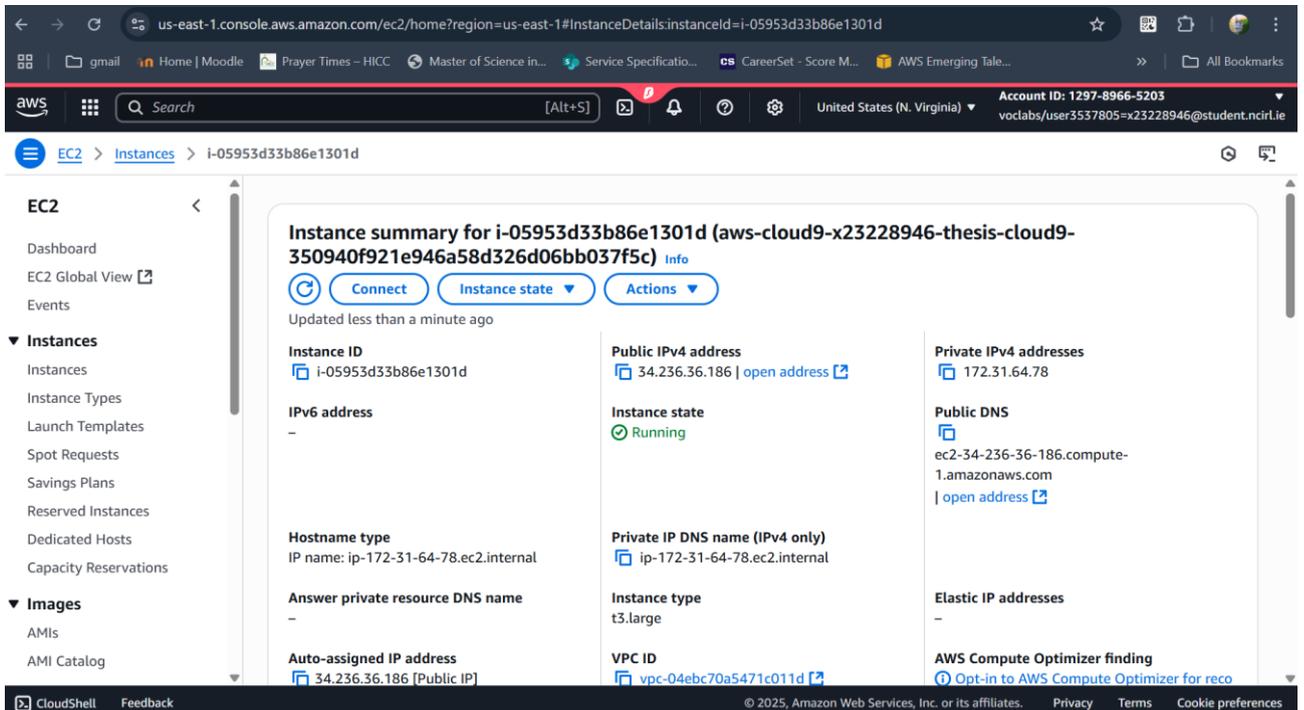
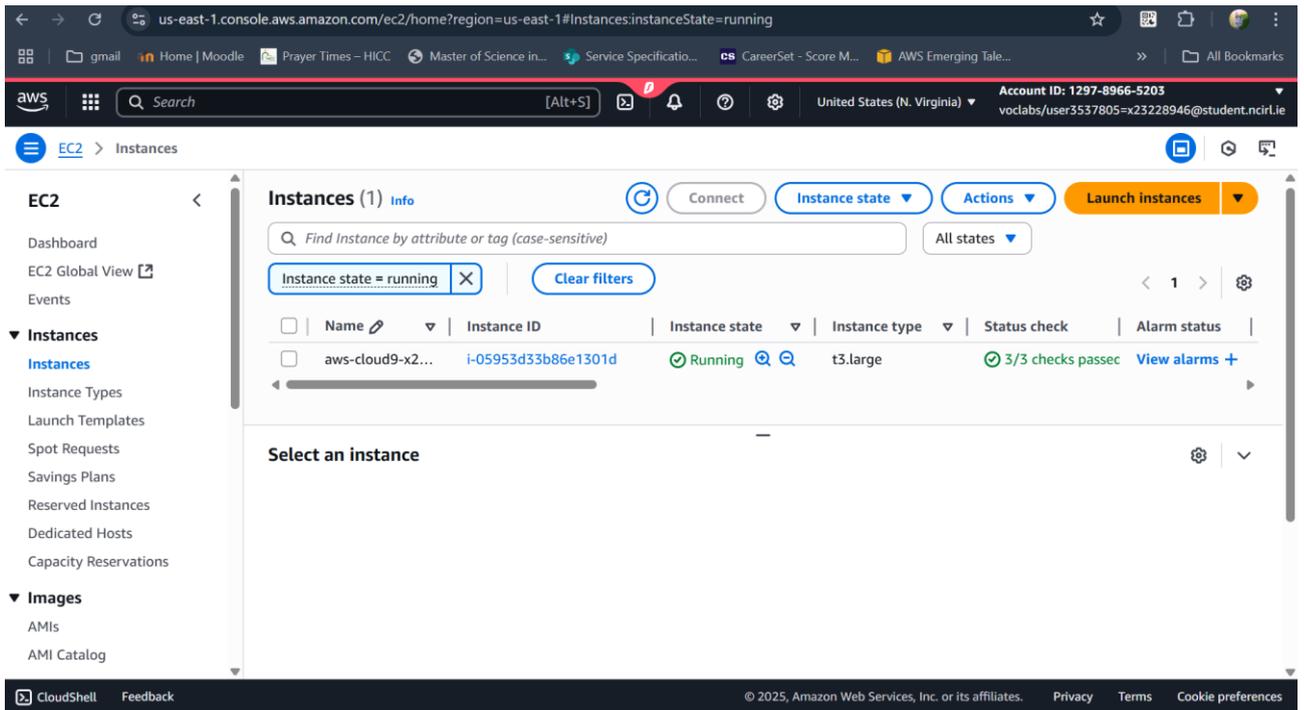


Figure 3: Data Loading Snippet

```

1]: #import all Library
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import os
import pickle
import plotly.graph_objects as go
import plotly.express as px
import plotly.figure_factory as ff
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.preprocessing import LabelEncoder, MinMaxScaler
from sklearn import metrics
from imblearn.under_sampling import NearMiss
from sklearn.metrics import precision_recall_fscore_support
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
import warnings
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
from lightgbm import LGBMClassifier
from xgboost import XGBClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import AdaBoostClassifier, GradientBoostingClassifier
from sklearn.tree import DecisionTreeClassifier
from mlxtend.classifier import StackingClassifier
# from pycaret.classification import *

```

Loading dataset from s3 bucket

```

[5]: # Config
bucket_name = 'x23228946-thesis-bucket'
file_key = 'Train_Test_Network.csv'

# Create S3 client
s3_client = boto3.client('s3')

# Download CSV file content to memory
response = s3_client.get_object(Bucket=bucket_name, Key=file_key)
csv_content = response['Body'].read().decode('utf-8')

# Load into DataFrame
iotdataset = pd.read_csv(StringIO(csv_content))

```

Random Forest

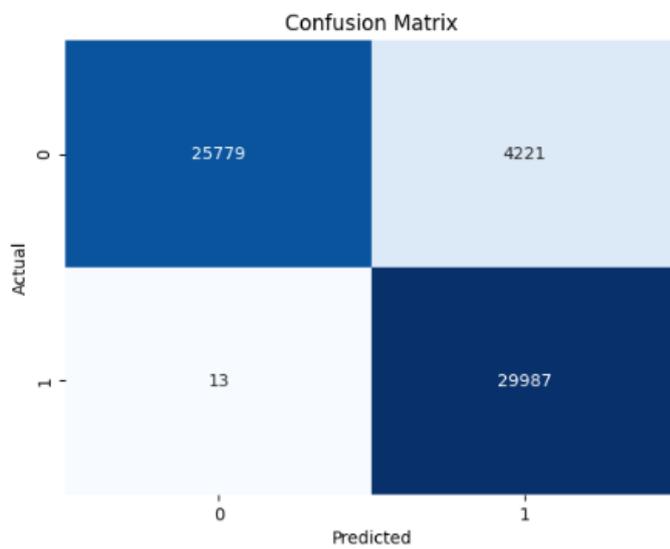
RANDOM FOREST

```
j]: rfmodel = RandomForestClassifier( max_depth=3,n_estimators=3)
rfmodel.fit(X_train,y_train)
y_pred = rfmodel.predict(X_test)

j]: matrix = confusion_matrix(y_test, y_pred)
sns.heatmap(matrix, annot=True, cbar=False, cmap='Blues', fmt="d")
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

# Save figure Locally
local_graph_path = '/home/ec2-user/environment/graphs/random_forest.png'
fig.write_image(local_graph_path)
s3_key = 'graphs/random_forest.png'
s3_client.upload_file(local_graph_path, bucket_name, s3_key)

print(f" Uploaded graph to s3://{bucket_name}/{s3_key}")
```



```
Classification Report :
      precision    recall  f1-score   support

     0       1.00      0.86      0.92     30000
     1       0.88      1.00      0.93     30000

 accuracy          0.94
 macro avg          0.93
 weighted avg       0.93
```

Gradient Boosting

GRADIENT BOOSTING

```
from sklearn.ensemble import AdaBoostClassifier, GradientBoostingClassifier
gb = GradientBoostingClassifier(max_depth=2, n_estimators=3)
gb.fit(X_train, y_train)
y_pred = gb.predict(X_test)
```

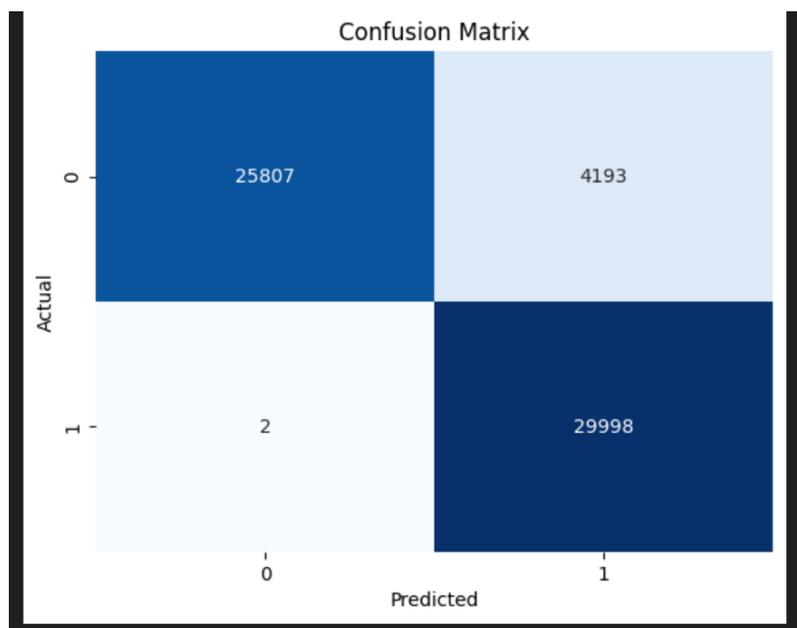
Python

```
matrix = confusion_matrix(y_test, y_pred)
sns.heatmap(matrix, annot=True, cbar=False, cmap='Blues', fmt="d")
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

# Save figure locally
local_graph_path = '/home/ec2-user/environment/graphs/gradient_boosting.png'
fig.write_image(local_graph_path)
s3_key = 'graphs/gradient_boosting.png'
s3_client.upload_file(local_graph_path, bucket_name, s3_key)

print(f" Uploaded graph to s3://{bucket_name}/{s3_key}")
```

Python



```
#Classification Report
print("Classification Report : ")
print(classification_report(y_test, y_pred))
```

```
Classification Report :
              precision    recall  f1-score   support

     0           1.00       0.86      0.92     30000
     1           0.88       1.00      0.93     30000

 accuracy                   0.93     60000
 macro avg                  0.94     60000
 weighted avg               0.94     60000
```

Naïve Bayes

NAIVE BAYES

```
from sklearn.naive_bayes import GaussianNB
nb=GaussianNB()
nb.fit(X_train,y_train)
y_pred = nb.predict(X_test)
```

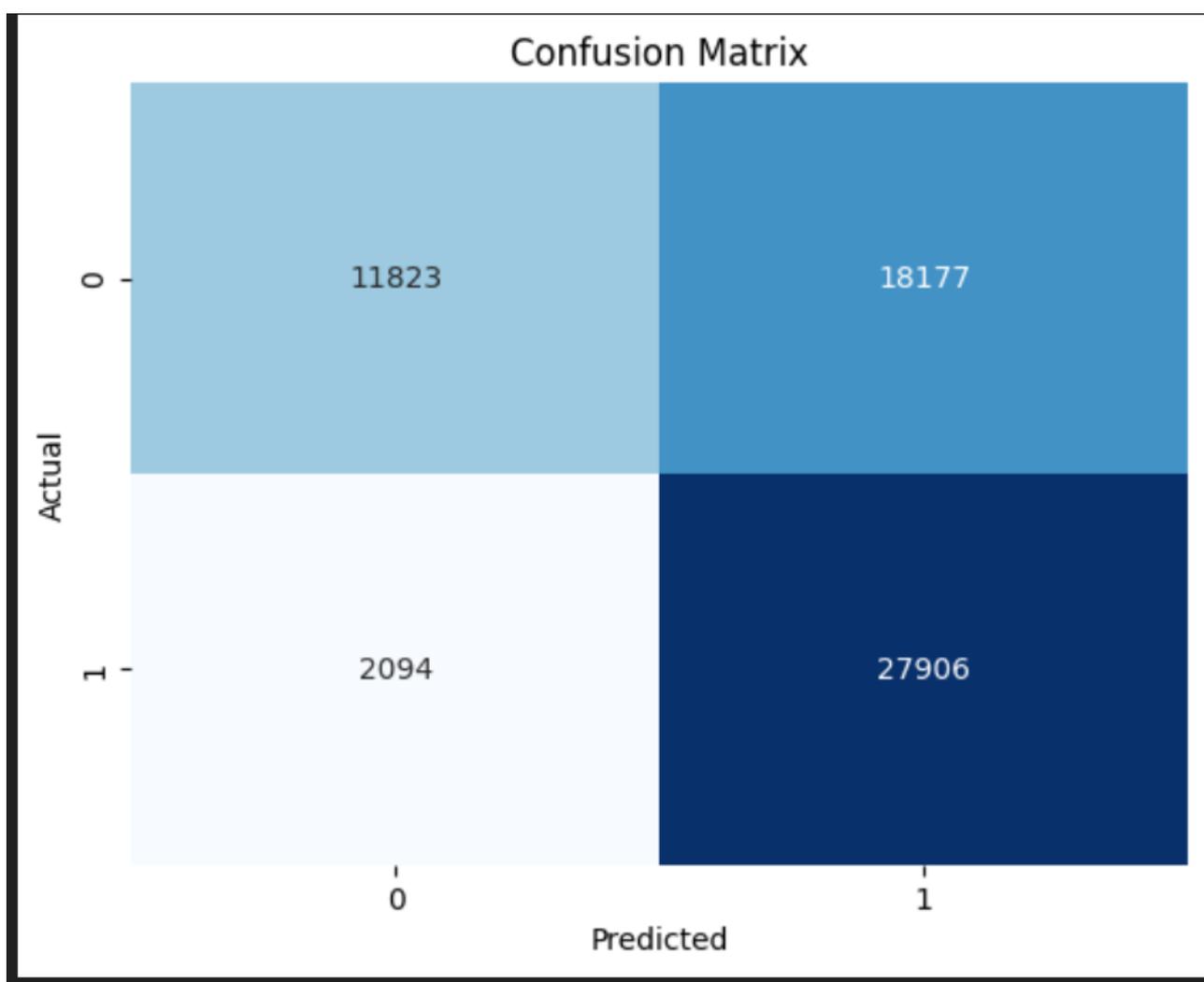
Python

```
matrix = confusion_matrix(y_test, y_pred)
sns.heatmap(matrix, annot=True, cbar=False, cmap='Blues', fmt="d")
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

# Save figure locally
local_graph_path = '/home/ec2-user/environment/graphs/naive_bayes.png'
fig.write_image(local_graph_path)
s3_key = 'graphs/naive_bayes.png'
s3_client.upload_file(local_graph_path, bucket_name, s3_key)

print(f" Uploaded graph to s3://{bucket_name}/{s3_key}")
```

Python



```
... Classification Report :
              precision    recall  f1-score   support

     0       0.85         0.39         0.54       30000
     1       0.61         0.93         0.73       30000

 accuracy         0.66       60000
 macro avg        0.73         0.66         0.64       60000
 weighted avg     0.73         0.66         0.64       60000
```

Logistic Regression:

LOGISTIC REGRESSION

```
from sklearn.linear_model import LogisticRegression
lasso_clf = LogisticRegression(penalty='l1', solver='liblinear', C=1.0)
lasso_clf.fit(X_train, y_train)
y_pred = lasso_clf.predict(X_test)
```

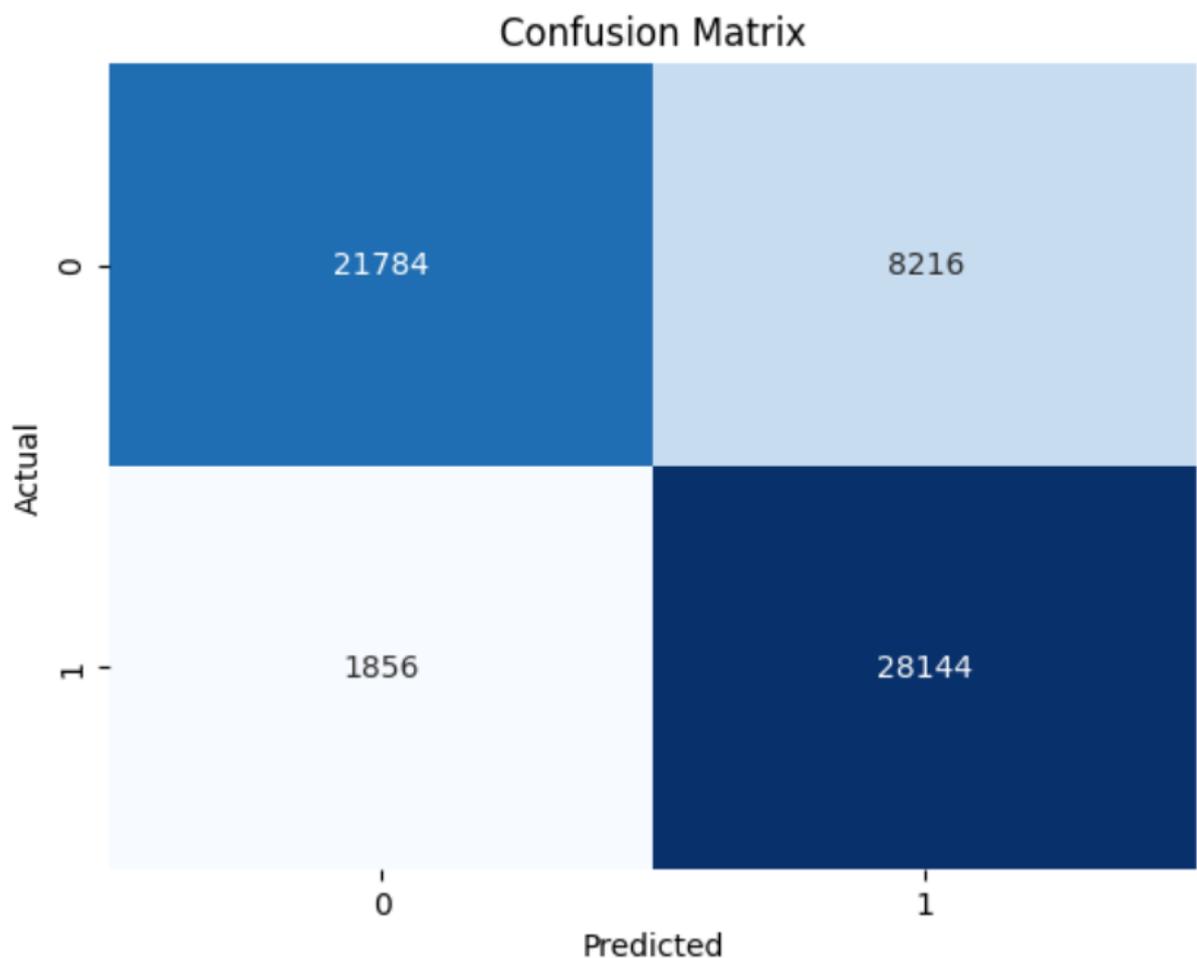
Python

```
matrix = confusion_matrix(y_test, y_pred)
sns.heatmap(matrix, annot=True, cbar=False, cmap='Blues', fmt="d")
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

# Save figure locally
local_graph_path = '/home/ec2-user/environment/graphs/logistic_regression.png'
fig.write_image(local_graph_path)
s3_key = 'graphs/logistic_regression.png'
s3_client.upload_file(local_graph_path, bucket_name, s3_key)

print(f" Uploaded graph to s3://{bucket_name}/{s3_key}")
```

Python



```

... Classification Report :
              precision    recall  f1-score   support

     0       0.92         0.73         0.81     30000
     1       0.77         0.94         0.85     30000

 accuracy                   0.83     60000
 macro avg                   0.85     60000
 weighted avg                0.85     60000

```

Superlearner:

```

def get_models():
    models = list()
    models.append(RandomForestClassifier( max_depth=3))
    models.append(LogisticRegression(penalty='l1', solver='liblinear', C=1.0))
    models.append(GaussianNB())
    return models

```

```

# collect out of fold predictions form k-fold cross validation
def get_out_of_fold_predictions(X, y, models):
    meta_X, meta_y = list(), list()
    # define split of data
    kfold = KFold(n_splits=5, shuffle=True)
    # enumerate splits
    for train_ix, test_ix in kfold.split(X):
        fold_yhats = list()
        # get data
        # Use .iloc to index by integer position
        train_X, test_X = X.iloc[train_ix], X.iloc[test_ix]
        # Assuming y is a numpy array or can be indexed directly by integer position
        train_y, test_y = y[train_ix], y[test_ix]
        meta_y.extend(test_y)
        # fit and make predictions with each sub-model
        for model in models:
            model.fit(train_X, train_y)
            yhat = model.predict(test_X)
            # store columns
            fold_yhats.append(yhat.reshape(len(yhat),1))
        # store fold yhats as columns
        meta_X.append(hstack(fold_yhats))
    return vstack(meta_X), asarray(meta_y)

```

```

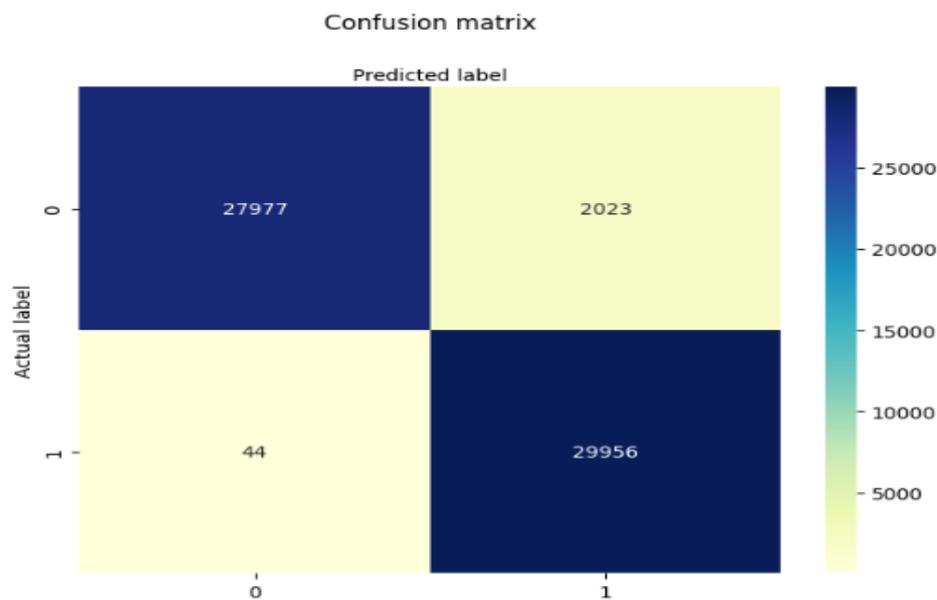
# fit all base models on the training dataset
def fit_base_models(X, y, models):
    for model in models:
        model.fit(X, y)

# fit a meta model
def fit_meta_model(X, y):
    model = GradientBoostingClassifier( max_depth=2,n_estimators=3)
    model.fit(X, y)
    return model

# evaluate a List of models on a dataset
def evaluate_models(X, y, models):
    for model in models:
        yhat = model.predict(X)
        acc = accuracy_score(y, yhat)
        print('%s: %.3f' % (model.__class__.__name__, acc*100))

# make predictions with stacked model
def super_learner_predictions(X, models, meta_model):
    meta_X = list()
    for model in models:
        yhat = model.predict(X)
        meta_X.append(yhat.reshape(len(yhat),1))
    meta_X = hstack(meta_X)
    # predict
    return meta_model.predict(meta_X)

```



Uploaded graph to s3://x23228946-thesis-bucket/graphs/super_learner.png

	precision	recall	f1-score	support
0	1.00	0.93	0.96	30000
1	0.94	1.00	0.97	30000
accuracy			0.97	60000
macro avg	0.97	0.97	0.97	60000
weighted avg	0.97	0.97	0.97	60000